# Pivotal.

# Agile Data Science on Greenplum Using Airflow

Ambarish Joshi - Data Scientist

Aditya Padhye - Data Engineer

"**Agile software development** refers to a group of **software development** methodologies based on **iterative development**, where requirements and solutions evolve through collaboration between self-organizing cross-functional teams."

Pivotal

# Data Science Phases

**Discovery Phase**

**Operationalization (O16n) Phase**

Pivotal

# Data Science Phases

**Discovery Phase**

✓ Data exploration & cleaning

✓ Feature engineering

✓ Model Building

✓ Model Evaluation

Pivotal

# Data Science Phases - Agility

**Discovery Phase**

✓ Data exploration & cleaning

✓ Feature engineering

✓ Model Building

✓ Model Evaluation

Rapid Iteration and Experimentation

Pivotal

# Data Science Phases - Agility

**Discovery Phase**

✓ Data exploration & cleaning

✓ Feature engineering

✓ Model Building

✓ Model Evaluation

jupyter

Pivotal
**Greenplum**

Rapid Iteration and Experimentation

Pivotal.

# Greenplum Database

- **MPP database based on Postgres**
- **In database analytics**
- **Parallel architecture**

# Jupyter Notebooks

# Data Science Phases

**Discovery Phase**

**Operationalization (O16n) Phase**

Pivotal

# Data Science Phases

**Operationalization (O16n) Phase**

✓ Data Pipelines

✓ Testing

✓ Monitoring

● APIs to consume model output

Pivotal.

# Data Science Phases - Agility



**Operationalization (O16n) Phase**

✓ Automated manageable pipelines

✓ Data Pipelines

✓ Testing with CI

✓ Testing

✓ Monitoring to react to Failures

✓ Monitoring

Madlib Flow Talk by Frank and Sridhar

APIs to consume model output

**Pivotal**

# Data Science Phases - Agility



### Operationalization (O16n) Phase

✓ Automated manageable pipelines

✓ Data Pipelines

✓ Testing with CI

✓ Testing

✓ Monitoring to react to Failures

✓ Monitoring

Madlib Flow Talk by Frank and Sridhar

APIs to consume model output

# Airflow

- **Apache Project spun out of Airbnb**

- **"Airflow is a platform to programmatically author, schedule and monitor workflows."**



Pivotal

# Data Science Use-Case

- **The Data**
  - **Time-series trajectories** with **latitude** and **longitude** of location.
  - Subset of trajectories are labeled as **walk** / **not walk**

- **Our Model**
  - Build **Classification model** using labelled data to identify if new unlabeled trajectories are walk or not walk

**Example trajectories**



Pivotal

# Example data

## Example trajectory data

| uid | latitude | longitude | tdate | ttime |
|-----|----------|-----------|-------|-------|
| 020 | 39.9744533333333 | 116.302163333333 | 2011-08-25 | 14:38:25 |
| 020 | 39.97445 | 116.302165 | 2011-08-25 | 14:38:26 |
| 020 | 39.9746016666667 | 116.302073333333 | 2011-08-25 | 14:39:22 |
| 020 | 39.97473 | 116.302066666667 | 2011-08-25 | 14:39:23 |
| 020 | 39.9748116666667 | 116.30207 | 2011-08-25 | 14:39:24 |

**We have mode labels of walk and not walk only for subset of incoming daily trajectories**

## Example label data

| uid | start_date | start_time | end_date | end_time | mode |
|-----|-----------|-----------|----------|----------|------|
| 020 | 2011-08-27 | 06:13:01 | 2011-08-27 | 08:01:37 | walk |
| 020 | 2011-08-27 | 09:34:43 | 2011-08-27 | 14:50:30 | walk |
| 020 | 2011-08-27 | 14:50:31 | 2011-08-27 | 15:01:58 | bus |
| 020 | 2011-08-27 | 15:01:59 | 2011-08-27 | 15:31:43 | walk |
| 020 | 2011-08-28 | 04:33:31 | 2011-08-28 | 04:44:25 | walk |

Pivotal

# Discovery phase ➜ Operationalization phase

**After every model iteration we check if the model is viable**

- Check the quantitative metrics of the model like AUC, ROC curve, accuracy etc
- Check the qualitative results of the model and if it make sense to a subject matter expert

**Once we are convinced that the model is both quantitatively and qualitatively viable we can move to the Operationalization phase**

Pivotal

# Discovery phase ➜ Operationalization phase

## Example of code from the discovery phase which is converted into a task script

# Architecture overview



Discovery Phase

Model Evaluation · Data Exploration · Iterative Discovery Phase · Feature Engineering · Model Building

RAW

output → ML Model + Data pre-processing

Operationalization Phase

Modular idempotent tasks

Connect tasks to create automated workflows

output → Model refitting workflow

output → New data inference workflow

TBD: Expose model results using an API

Pivotal

# Data Science Phases - Agility

**Pivotal Greenplum**

**Operationalization (O16n) Phase**

✓ Automated manageable Pipelines

✓ Pipelines

✓ Testing with CI/CD

✓ Testing

✓ Monitoring to React to Failures

✓ Monitoring

Madlib Flow Talk by Frank and Sridhar

APIs to consume model output

**Pivotal**

# Data Prep and Feature Engineering



| | | |
|---|---|---|
| **Fetch** | fetch_daily_trajectory | fetch_daily_label |
| **Clean** | clean_daily_trajectory | clean_daily_label |
| **Transform** | merge_trajectory_label | |
| | calculate_trajectory_speed | |
| | trajectory_speed_walk | |
| **Feature Engineering** | create_tsfresh_features | |
| | pivot_tsfresh_features | |
| Extract **labelled data** for model **creation/refitting** | tsfresh_model_features | tsfresh_predict_features |
| **Inference** for **unlabelled data** | | predict_walk_trajectories |

Pivotal

# Data Prep and Feature Engineering - Demo

Pivotal

# Model Training



| | | |
|---|---|---|
| **Fetch** | fetch_daily_trajectory | fetch_daily_label |
| **Clean** | clean_daily_trajectory | clean_daily_label |
| **Transform** | merge_trajectory_label | |
| | calculate_trajectory_speed | |
| | trajectory_speed_walk | |
| **Feature Engineering** | create_tsfresh_features | |
| | pivot_tsfresh_features | |
| Extract **labelled data** for model **creation/refitting** | tsfresh_model_features | tsfresh_predict_features |
| **Inference** for **unlabelled data** | | predict_walk_trajectories |

Pivotal.

# Model Training

- This DAG has a single task for model training
- In this task we split the data into train and test samples, train the model, evaluate the model and capture the accuracy, auc and model tables.
- We want all of the above to run at the same time

```
train_model
```

```sql
SELECT madlib.train_test_split(
                    'geolife.tsfresh_model_features',     -- Source table
                    'geolife.features_walk_{{ds_nodash}}',     -- Output table
                    0.8,          -- Sample proportion
                    0.2,          -- Sample proportion
                    NULL, -- Strata definition
                    NULL, -- Columns to output
                    FALSE,       -- Sample without replacement
                    TRUE);      -- Separate output tables

-- build a random forest model using madlib

DROP TABLE IF EXISTS geolife.rf_walk_{{ds_nodash}}_output, geolife.rf_walk_{{ds_nodash}}_outp
SELECT madlib.forest_train('geolife.features_walk_{{ds_nodash}}_train',        -- source tab
                    'geolife.rf_walk_{{ds_nodash}}_output',       -- output model table
                    'id',          -- id column
                    'label',          -- response
                    '*',   -- features
                    'tdate',            -- exclude columns
                    NULL,           -- grouping columns
                    20::integer,       -- number of trees
                    2::integer,        -- number of random features
                    TRUE::boolean,      -- variable importance
                    1::integer,        -- num_permutations
                    8::integer,        -- max depth
                    3::integer,        -- min split
                    1::integer,        -- min bucket
                    10::integer       -- number of splits per continuous variable
                    );

-- Evaluate the built model

DROP TABLE IF EXISTS geolife.rf_walk_{{ds_nodash}}_results;
SELECT madlib.forest_predict('geolife.rf_walk_{{ds_nodash}}_output',       -- tree model
                    'geolife.features_walk_{{ds_nodash}}_test',        -- new d
                    'geolife.rf_walk_{{ds_nodash}}_results') ;--, -- output table
                    --'prob');

-- Capture model results

drop table if exists geolife.walk_{{ds_nodash}}_result;
create table geolife.walk_{{ds_nodash}}_result
as
with t as (
select id,
    case when label = True then 1.0 else 0.0 end as obs
from geolife.features_walk_{{ds_nodash}}_test
```

Pivotal

# Model Training - Demo

Pivotal

# Model Scoring

# Model Scoring

- The unlabeled data which is extracted from the features table is scored in this DAG
- We first check if any model has been built
- If there is a model so we score the data (inference)

```
DO
$do$
DECLARE
tabname    character varying(255);
BEGIN
IF (select count(*) from
  (select 1 from geolife.tsfresh_predict_features{{ds_nodash}}  limit 1) as t) > 0
  and
  (select count(*) from (select 1 from geolife.models_metadata  limit 1) as p) > 0
THEN
    tabname := (select model_tabname
    from geolife.models_metadata
    order by mdate
    limit 1);

    DROP TABLE IF EXISTS prediction_results;
    PERFORM madlib.forest_predict(tabname,
                            'geolife.tsfresh_predict_features{{ds_nodash}}',
                            'geolife.walk_prediction_results{{ds_nodash}}',
                            'response');

    insert into geolife.walk_prediction_results
    select *,
        regexp_replace(id, '^.*([0-9-]{10})_.*$', E'\\1')::date as tdate
    from geolife.walk_prediction_results{{ds_nodash}};

END IF;
END
$do$;

drop table if exists geolife.tsfresh_predict_features{{ds_nodash}};
drop table if exists geolife.walk_prediction_results{{ds_nodash}};
```

Pivotal.

Model Scoring - Demo

# Model Re-Training

- Daily we get some more labeled data, once we have accumulated enough labeled data we can retrain the model for better accuracy
- We have scheduled model re-training monthly

Pivotal.

Model Re-Training - Demo

# Data Science Phases - Agility

### Pivotal Greenplum

- ✓ Automated manageable Pipelines

- ✓ Testing with CI/CD

- ✓ Monitoring to React to Failures

  Madlib Flow Talk by Frank and Jarrod

### Operationalization (O16n) Phase

- ✓ Pipelines

- ✓ Testing

- ✓ Monitoring

  APIs to consume model output

Pivotal.

# Testing with CI/CD

- **Testing Data Pipelines is hard**

- **Test Coverage (Test Tasks vs Test DAGs)**

- **Testing as part of the CI/CD**

Pivotal

# Testing with CI/CD - Demo

# Data Science Phases - Agility

Pivotal **Greenplum**

## Operationalization (O16n) Phase

✓ Automated manageable Pipelines

✓ Testing with CI/CD

✓ Monitoring to React to Failures

Madlib Flow Talk by Frank and Sridhar

✓ Pipelines

✓ Testing

✓ Monitoring

APIs to consume model output

**Pivotal.**

# Monitoring and Error Fixing

- Monitoring and error fixing is big part of responsive data pipelines

- Ability to quickly identify what is failing, why it is failing and fixing it with minimum lead time is crucial

- In this demo we will showcase an error fixing case

Pivotal

# Monitoring and Error Fixing - Demo

Pivotal

# Data Science Phases - Agility

**Pivotal Greenplum**

**Operationalization (O16n) Phase**

✓ Automated manageable Pipelines

✓ Testing with CI/CD

✓ Monitoring to React to Failures

Madlib Flow Talk by Frank and Sridhar

✓ Pipelines

✓ Testing

✓ Monitoring

APIs to consume model output

**Pivotal**

# Conclusion

✓ **Greenplum and Jupyter notebooks provides a set of tools to do Agile Data Science during discovery phase**

✓ **Greenplum along with Airflow and Circle CI is very effective to do Agile Data Science during the operationalization phase**

Pivotal

# Questions

Pivotal